

Stochastic Differential Equations

Guanyu Chen

January 27, 2025

Abstract

This paper is about the basic theory and applications of SDEs. Also, this is a sum up of last term's SDE course. Reference: [6][1][5] [3][4][7][2]

Contents

1	Introduction to SDEs	3
1.1	SODEs	3
1.2	The It'so and Stratonovich Stochastic Integrals	3
1.3	Ito's Formula	3
1.4	Mean and Covariance	3
2	Fokker-Planck-Kolmogorov Equation	4
2.1	FPK Equation	4
2.2	Langevin SDE	6
3	Numerical Methods	7
4	What is Diffusion After All?	7
4.1	From SDEs	7
4.2	From Flow Map	8
4.3	Solution	8
5	Variational Auto-Encoder	9
5.1	Structure	9
5.2	Evidence Lower Bound	10
6	Diffusion Model	10
6.1	DDPM	10
6.2	Score Matching	13
6.2.1	Explicit Score Matching	14
6.2.2	Implicit Score Matching	14
6.2.3	Denoising Score Matching	14
6.3	Denoising Diffusion	15
6.3.1	With Classifier Guidance	15
6.3.2	Classifier Guidance Free	15
7	Flow Matching	16
7.1	FPK Equation	16
7.2	Probability Flow	16
8	Semilinear Stochastic PDEs	17
8.1	Semilinear SPDE	17
8.2	Q wiener process	17
8.3	Cylindrical Wiener Process	18
8.4	Ito integral solution	18
8.5	Semilinear SPDE	18
A	Conservation Laws	21
A.1	Flow Map	21
A.1.1	Conservation Laws	21

B	Priori	22
B.1	Hilbert space-valued random variable	22
B.2	Hilbert-Schmidt operator	22
B.3	Operator theory	22

1 Introduction to SDEs

1.1 SODEs

Problem 1 Assume we have a Stochastic Differential Equation like:

$$dX_t = f(X_t, t)dt + G(X_t, t)dW_t \quad (1)$$

where $X_t \in \mathbf{R}^d$, $f \in \mathcal{L}(\mathbf{R}^{d+1}, \mathbf{R}^d)$, and W_t is m -dim Brownian Motion with diffusion matrix Q , $G(X_t, t) \in \mathcal{L}(\mathbf{R}^{m+1}, \mathbf{R}^d)$, with initial condition $X_0 \sim p(X_0)$.

1.2 The It'so and Stratonovich Stochastic Integrals

1.3 Ito's Formula

1.4 Mean and Covariance

We can derive the mean and covariance of SDE. By applying Ito's formula to $\phi(x, t)$, then

$$\frac{dE[\phi]}{dt} = E\left[\frac{\partial \phi}{\partial t}\right] + \sum_i E\left[\frac{\partial \phi}{\partial x_i} f_i(X_t, t)\right] + \frac{1}{2} \sum_{ij} E\left[\frac{\partial^2 \phi}{\partial x_i \partial x_j} [GG^T]_{ij}\right] \quad (2)$$

By taking $\phi(X, t) = x_i$ and $\phi(X, t) = x_i x_j - m(t)_i m(t)_j$, we have the mean function $m(t) = E[X_t]$ and covariance function $c(t) = E[(X_t - m(t))(X_t - m(t))^T]$ respectively, s.t.

$$\begin{cases} \frac{dm}{dt} = E[f(X_t, t)] \\ \frac{dc}{dt} = E[f(X, t)(X - m(t))^T] + E[(X - m(t))f^T(X, t)] + E[G(X_t, t)QG^T(X_t, t)] \end{cases} \quad (3)$$

So we can estimate the mean and covariance of solution to SDE. However, these equations cannot be used as such, because only in the Gaussian case do the expectation and covariance actually characterize the distribution.

The linear SDE has explicit solution. Assume the linear SDE has the form

$$dX_t = (K(t)X_t + B(t))dt + G(t)dW_t \quad (4)$$

where $K(t) \in \mathbf{R}^{d \times d}$, $B(t) \in \mathbf{R}^d$, $G(t) \in \mathbf{R}^{d \times m}$ are given functions. $X_t \in \mathbf{R}^d$ is the state vector, $W_t \in \mathbf{R}^m$ is the Brownian Motion with diffusion matrix Q .

Theorem 1 The explicit solution to the linear SDE is given by:

$$X_t = \Psi(t, t_0)X_0 + \int_{t_0}^t \Psi(t, s)B(s)ds + \int_{t_0}^t \Psi(t, s)G(s)dW_s \quad (5)$$

where $\Psi(t, t_0)$ is the transition matrix of the linear SDE, which satisfies the following matrix ODE:

$$\frac{d\Psi}{dt} = K(t)\Psi(t, t_0), \Psi(t_0, t_0) = I \quad (6)$$

Hence, X_t is a Gaussian process (A linear transformation of Brownian Motion which is a Gaussian process).

Proof 1 Multiply both sides of the SDE by Integrating factor $\Psi(t_0, t)$ and apply Ito's formula to $\Psi(t_0, t)X_t$. See Sarkka P49.

As discussed above, we can compute the mean and covariance function of solution to linear SDE.

Theorem 2 The mean and covariance function of solution to linear SDE are given by:

$$\begin{cases} \frac{dm}{dt} = K(t)m(t) + B(t) \\ \frac{dc}{dt} = K(t)c(t) + c(t)K^T(t) + G(t)QG^T(t) \end{cases} \quad (7)$$

with initial condition $m_0 = m(t_0) = E[X_0]$, $c_0 = c(t_0) = \text{Cov}(X_0)$. Then the solution is given by solving the above ODEs:

$$\begin{cases} m(t) = \Psi(t, t_0)m_0 + \int_{t_0}^t \Psi(t, s)B(s)ds \\ c(t) = \Psi(t, t_0)c_0\Psi^T(t, t_0) + \int_{t_0}^t \Psi(t, s)G(s)QG^T(s)\Psi^T(t, s)ds \end{cases} \quad (8)$$

Proof 2 Apply $F(X, t) = K(t)X + B(t)$, $G(X, t) = G(t)$ to 3.

Hence the solution to linear SDE is a Gaussian process with mean and covariance function given by the above ODEs.

Theorem 3 The solution to LSDE is Gaussian:

$$p(X, t) = \mathcal{N}(X(t)|m(t), c(t)) \quad (9)$$

Specially when $X_0 = x_0$ is fixed, then

$$p(X, t|X_0 = x_0) = \mathcal{N}(X(t)|m(t|x_0), c(t|x_0)) \quad (10)$$

That is, $m_0 = x_0, c_0 = 0$. Then we have:

$$\begin{cases} m(t|x_0) = \Psi(t, t_0)x_0 + \int_{t_0}^t \Psi(t, s)B(s)ds \\ c(t|x_0) = \int_{t_0}^t \Psi(t, s)G(s)QG^T(s)\Psi^T(t, s)ds \end{cases} \quad (11)$$

Proof 3 The proof is straight forward either by applying $m_0 = x_0, c_0 = 0$ to 8 or by eq 5.

So, to sum up, linear SDE has great properties! The distribution is completely decided by the initial condition. Also, if we generate X_0 to X_{t_k} , which means that we begin SDE at t_i with X_{t_i} , we have the equivalent discretization of SDE:

Theorem 4 Original SDE is weakly, in distribution, equivalent to the following discrete-time SDE:

$$X_{t_{i+1}} = A_i X_{t_i} + B_i + G_i \quad (12)$$

where

$$\begin{cases} A_i = \Psi(t_{i+1}, t_i) \\ B_i = \int_{t_i}^{t_{i+1}} \Psi(t_{i+1}, s)B(s)ds \\ G_i = \int_{t_i}^{t_{i+1}} \Psi(t_{i+1}, s)G(s)QG^T(s)\Psi^T(t_{i+1}, s)ds \end{cases} \quad (13)$$

Proof 4 The proof is straight forward.

Theorem 5 The covariance of X_t and $X_s (s < t)$ is given by:

$$\text{Cov}(X_t, X_s) = \Psi(t, s)c(s) \quad (14)$$

Proof 5 See Sarkka P88-89.

2 Fokker-Planck-Kolmogorov Equation

2.1 FPK Equation

Definition 1 (Generator) The infinitesimal generator of a stochastic process $X(t)$ for function $\phi(x)$, i.e. $\phi(X_t)$ can be defined as

$$\mathcal{A}\phi(X_t) = \lim_{s \rightarrow 0^+} \frac{E[\phi(X(t+s)) - \phi(X(t))]}{s} \quad (15)$$

Where ϕ is a suitable regular function.

This leads to Dynkin's Formula very naturally.

Theorem 6 (Dynkin's Formula)

$$E[f(X_t)] = f(X_0) + E\left[\int_0^t \mathcal{A}(f(X_s))ds\right] \quad (16)$$

Theorem 7 If $X(t)$ s.t. [1](#), then the generator is given:

$$\mathcal{A}(\cdot) = \sum_i \frac{\partial(\cdot)}{\partial x_i} f_i(X_t, t) + \frac{1}{2} \sum_{i,j} \left(\frac{\partial^2(\cdot)}{\partial x_i \partial x_j} \right) [G(X_t, t) Q G^\top(X_t, t)]_{ij} \quad (17)$$

Proof 6 See P119 of SDE by Oksendal.

Example 1 If $dX_t = dW_t$, then $\mathcal{A} = \frac{1}{2}\Delta$, where Δ is the Laplace operator.

Definition 2 (Generalized Generator) For $\phi(x, t)$, i.e. $\phi(X_t, t)$, the generator can be defined as:

$$A_t \phi(x, t) = \lim_{s \rightarrow 0^+} \frac{E[\phi(X(t+s), t+s)] - \phi(X(t), t)}{s} \quad (18)$$

Theorem 8 Similarly if $X(t)$ s.t. [1](#), then the generalized generator is given:

$$\mathcal{A}_t(\cdot) = \frac{\partial(\cdot)}{\partial t} + \sum_i \frac{\partial(\cdot)}{\partial x_i} f_i(X_t, t) + \frac{1}{2} \sum_{i,j} \left(\frac{\partial^2(\cdot)}{\partial x_i \partial x_j} \right) [G(X_t, t) Q G^\top(X_t, t)]_{ij} \quad (19)$$

We want to consider the density distribution of $X_t, P(x, t)$

Theorem 9 (Fokken-Planck-Kolmogorov equation) The density function $P(x, t)$ of X_t s.t. [1](#) solves the PDE:

$$\frac{\partial P(x, t)}{\partial t} = - \sum_i \frac{\partial}{\partial x_i} [f_i(x, t) p(x, t)] + \frac{1}{2} \sum_{i,j} \frac{\partial^2}{\partial x_i \partial x_j} [(G Q G^\top)_{ij} P(x, t)] \quad (20)$$

The PDE is called FPK equation / forward Kolmogorov equation.

Proof 7 Consider the function $\phi(x)$, let $x = X_t$ and apply Ito's Formula:

$$\begin{aligned} d\phi &= \sum_i \frac{\partial \phi}{\partial x_i} dx_i + \frac{1}{2} \sum_{i,j} \left(\frac{\partial^2 \phi}{\partial x_i \partial x_j} \right) dx_i dx_j \\ &= \sum_i \frac{\partial \phi}{\partial x_i} (f_i(X_t, t) dt + (G(X_t, t) dW_t)) + \frac{1}{2} \sum_{i,j} \left(\frac{\partial^2 \phi}{\partial x_i \partial x_j} \right) [G(X_t, t) Q G^\top(X_t, t)]_{ij} dt. \end{aligned} \quad (21)$$

Take expectation of both sides:

$$\frac{dE[\phi]}{dt} = \sum_i E \left[\frac{\partial \phi}{\partial x_i} f_i(X_t, t) \right] + \frac{1}{2} \sum_{ij} E \left[\frac{\partial^2 \phi}{\partial x_i \partial x_j} [G Q G^\top]_{ij} \right] \quad (22)$$

So

$$\begin{cases} \frac{dE[\phi]}{dt} = \frac{d}{dt} \left[\int \phi(x) P(X_t = x, t) dx \right] = \int \phi(x) \frac{\partial P(x, t)}{\partial t} dx \\ \sum_i E \left[\frac{\partial \phi}{\partial x_i} f_i \right] = \sum_i \int \frac{\partial \phi}{\partial x_i} f_i(X_t = x, t) P dx = - \sum_i \int \phi \cdot \frac{\partial}{\partial x_i} [f_i(x, t) p(x, t)] dx. \\ \frac{1}{2} \sum_{ij} E \left[\frac{\partial^2 \phi}{\partial x_i \partial x_j} [G Q G^\top]_{ij} \right] = \frac{1}{2} \sum_{ij} \int \frac{\partial^2 \phi}{\partial x_i \partial x_j} [G Q G^\top]_{ij} P dx = \frac{1}{2} \sum_{ij} \int \phi(x) \frac{\partial^2}{\partial x_i \partial x_j} ([G Q G^\top]_{ij} P) dx. \end{cases} \quad (23)$$

then

$$\int \phi \frac{\partial P}{\partial t} dX = - \sum_i \int \phi \frac{\partial}{\partial x_i} (f_i P) dX + \frac{1}{2} \sum_{ij} \int \phi \frac{\partial^2}{\partial x_i \partial x_j} ([G Q G^\top]_{ij} P) dX$$

Hence

$$\int \phi \cdot \left[\frac{\partial P}{\partial t} + \sum_i \frac{\partial}{\partial x_i} (f_i P) - \frac{1}{2} \sum_{ij} \frac{\partial^2}{\partial x_i \partial x_j} ([G Q G^\top]_{ij} P) \right] dX = 0$$

Therefore P s.t.

$$\frac{\partial P}{\partial t} + \sum_i \frac{\partial}{\partial x_i} (f_i(x, t) P(x, t)) - \frac{1}{2} \sum_{i=1} \frac{\partial^2}{\partial X_i \partial X_j} ([G Q G^\top]_{ij} P(x, t)) = 0 \quad (24)$$

Which gives the FPK Equation.

Remark 1 When SDE is time independent:

$$dX_t = f(X_t)dt + G(X_t)dW_t \quad (25)$$

then the solution of FPK often converges to a stationary solution s.t. $\frac{\partial P}{\partial t} = 0$.

Here is an another way to show FPK equation: Since we have inner product $\langle \phi, \psi \rangle = \int \phi(x)\psi(x)dx$. Then $E[\phi(x)] = \langle \phi, P \rangle$.

As the equation 22 can be written as

$$\frac{d}{dt} \langle \phi, P \rangle = \langle \mathcal{A}\phi, P \rangle \quad (26)$$

Where \mathcal{A} has been mentioned above. If we note the adjoint operator of \mathcal{A} as \mathcal{A}^* , then we have

$$\langle \phi, \frac{dP}{dt} - \mathcal{A}^*(P) \rangle = 0, \forall \phi(x) \quad (27)$$

Hence we have

Theorem 10 (FPK Equation)

$$\frac{dP}{dt} = \mathcal{A}^*(P), \text{ where } \mathcal{A}^*(\cdot) = - \sum_i \frac{\partial}{\partial x_i} (f_i(x, t)(\cdot)) + \frac{1}{2} \sum_{i,j} \frac{\partial^2}{\partial x_i \partial x_j} ([GQG^\top]_{ij}(\cdot)) \quad (28)$$

It can be rewritten as:

$$\begin{aligned} \frac{\partial P}{\partial t} &= -\nabla \cdot [f(x, t)p(x, t)] + \frac{1}{2} \nabla^2 \cdot [(GQG^\top) p(x, t)] \\ &= -\nabla \cdot \left[f(x, t)p(x, t) - \frac{1}{2} \nabla \cdot [(GQG^\top) p(x, t)] \right] \end{aligned} \quad (29)$$

Theorem 11 (Transition Density(Forward Komogorov Equation)) The transition density $P_{t|s}(x_t|x_s), t \geq s$, which means the propability of transition from $X(s) = x_s$ to $X(t) = x_t$, satisfies the FPK equation with initial condition $P_{s|s}(x|x_s) = \delta(x - x_s)$ i.e. for $P_{t|s}(x|y)$, it solves

$$\frac{\partial P_{t|s}(x|y)}{\partial t} = \mathcal{A}^*(P_{t|s}(x|y)), \text{ with } P_{s|s}(x|y) = \delta(x - y) \quad (30)$$

Theorem 12 (Backward Komogorov Equation) $P_{s|t}(y|x)$ for $t \geq s$ solves:

$$\frac{\partial P_{s|t}(y|x)}{\partial s} + \mathcal{A}(P_{s|t}(y|x)) = 0, \text{ with } P_{s|s}(y|x) = \delta(x - y) \quad (31)$$

2.2 Langevin SDE

The Langevin SDE has the following form:

$$X_{t+s} = X_t + \nabla \log p_t(x_t)s + \sqrt{2s}\xi \quad (32)$$

where $X_t \in \mathcal{R}^d, p_t(x_t) = p(X_t = x_t), \xi \sim N(0, I), I$ is identical matrix of $m \times m$. Our goal is to sample from specific $p(x, t)$.

Theorem 13 The density of Langevin Diffusion Model converges to $p(x)$ over time. In other words, if $X_t \sim p(x)$, then $X_{t+s} \sim p(x)$ for $\forall s > 0$.

Proof 8 Let $\mu_t(f) = E[f(X_t)]$. Consider $\mu_{t+\tau}(f) = E[f(X_{t+\tau})]$, as $\tau \rightarrow 0$. Then

$$\begin{aligned} \mu_{t+\tau} &= E \left[f \left(X_t + \nabla \log p_t(x_t) \cdot \tau + \sqrt{2\tau}\xi \right) \right] \\ &= E \left[f(x_t) + \nabla^\top f(x_t) \left(\tau \nabla \log p_t(x_t) + \sqrt{2\tau}\xi \right) \right. \\ &\quad \left. + \frac{1}{2} \left(\nabla^\top \log p_t(x_t) \tau + \sqrt{2\tau}\xi \right) \nabla^2 f(x_t) \nabla \log p_t(x_t) \tau + \sqrt{2\tau}\xi \right] \\ &= E[f(x_t)] + E[\tau \nabla^\top f(x_t) \nabla \log p_t(x_t)] \\ &\quad + \frac{\tau^2}{2} E[\nabla^\top \log p(x_t) \cdot \nabla^2 f(x_t) \cdot \nabla \log p(x_t)] + E[\tau \xi^\top \nabla^2 f(x_t) \xi] \end{aligned} \quad (33)$$

The second term:

$$\begin{aligned}
& \tau E [\nabla^\top f \nabla \log p_t] \\
&= \tau \int \nabla f \cdot \nabla \log p_t p_t dx = \tau \int \nabla f \cdot \nabla p_t dx \\
&= -\tau \int \text{tr} (\nabla^2 f) \cdot p_t dx = -\tau E [\text{tr} (\nabla^2 f)] \\
&= -\tau E [\xi^\top \nabla^2 f \xi]
\end{aligned} \tag{34}$$

Then

$$\mu_{t+\tau} = E \left[\frac{1}{2} \nabla^\top \log p_t \nabla^2 f \nabla \log p_t \right] \cdot \tau^2 = O(\tau^2) \tag{35}$$

Hence we have $\frac{d}{dt}(\mu_t) = 0$, i.e. $E[\mu_t] = E[\mu_{t+s}]$ for $\forall s > 0$.

Remark 2 We define the density of normal distribution $N(x; \mu, \Sigma)$, and its log-density, gradient of density and score as follows:

$$\begin{cases} N(x; \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} e^{-\frac{1}{2}(x-\mu)^\top \Sigma^{-1}(x-\mu)} \\ \log N(x; \mu, \Sigma) = -\frac{1}{2}(x-\mu)^\top \Sigma^{-1}(x-\mu) - \log \left(\sqrt{(2\pi)^d |\Sigma|} \right) \\ \nabla_x N(x; \mu, \Sigma) = N(x; \mu, \Sigma) \Sigma^{-1}(x-\mu) \\ \nabla_x \log N(x; \mu, \Sigma) = -\Sigma^{-1}(x-\mu). \end{cases} \tag{36}$$

Actually, Langevin SDE is not necessary be as above i.e. the diffusion term is not necessary to be $\sqrt{2}$. The reason is to guarantee the stationary distribution of $p_t(x)$. i.e. the term $\frac{\partial p(x,t)}{\partial t} = 0$ in FPK equation. If the diffusion term is $g(t)$, then by FPK equation, we have

$$\nabla_x \cdot (fp - \frac{1}{2}g^2(t)\nabla p) = 0$$

then $f(x, t) = \frac{1}{2}g^2(t) \frac{\nabla_x p(x,t)}{p(x,t)} = \frac{1}{2}g^2(t)\nabla_x \log p(x, t)$.

3 Numerical Methods

Euler-Maruyama and Milstein methods have been talked before.

4 What is Diffusion After All?

4.1 From SDEs

At the beginning, the diffusion phenomenon is observed through the motion of particles(Brownian motion). Normally, the SDE can be written as:

$$dX_t = f(X_t, t)dt + G(X_t, t)dW_t \tag{37}$$

Here, we skip the drift term $f(X_t, t)$ and only consider the diffusion term $G(X_t, t)dW_t$, i.e.

$$dX_t = G(X_t, t)dW_t \tag{38}$$

Then by FPK equation, we can derive

Theorem 14 The probability density function $p(x, t)$ satisfies:

$$\frac{\partial p(x, t)}{\partial t} = \frac{1}{2} \sum_{i,j} \frac{\partial^2}{\partial x_i \partial x_j} \left[(GQG^\top)_{ij} p(x, t) \right] = \frac{1}{2} \nabla \cdot (\nabla \cdot (GQG^\top p(x, t))) \tag{39}$$

Specially, when $G(X_t, t) = G(t)$ and $Q = I$, we have:

$$\frac{\partial p}{\partial t} = \nabla \cdot \left(\frac{GG^\top}{2} \nabla p \right) \tag{40}$$

So, when $X_0 \sim p_0$, we can then compute the diffusion density $p(x, t)$ by solving the FPK equation.

4.2 From Flow Map

Since we have the definition of **Flow Map** $\phi_s^t(\mathbf{x})$, which is controlled by vector field $V(\phi_s^t(\mathbf{x}), t)$, then just think the $\phi_0^t(\mathbf{x})$ as the trajectory of the particle beginning at x over time, noted as $\phi_t(x)$. Then the vector field is actually the velocity field of the particle, so we have:

$$\begin{cases} \frac{\partial \phi_t(\mathbf{x})}{\partial t} = V(\phi_t(\mathbf{x}), t) \\ \phi_0(\mathbf{x}) = \mathbf{x} \end{cases} \quad (41)$$

The motion of particles described by ϕ_t determines how the density $p_t(x)$ evolves over time.

Theorem 15 *When the initial density $p_0(x)$ is known, the density field can be expressed as:*

$$p(\phi_t(x), t) = \frac{p_0(x)}{|\det J_{\phi_t}(x)|} \quad (42)$$

It should be noted that $\phi_t(x)$ is actually the same as X_t in SDE, then similarly, the density is:

$$\phi_t(x) \sim p_t(x) \quad (43)$$

So, the flow map is an ODE, which is a special case of SDE without diffusion term. Then we have:

Theorem 16 (Continuity Equation) *The probability density function $p(x, t)$ of X_t satisfies:*

$$\frac{\partial p(x, t)}{\partial t} = -\nabla \cdot (V(x, t)p(x, t)) \quad (44)$$

which is called **Continuity Equation**.

Remark 3 *The continuity equation can also be derived from the Conservation of Mass.*

Theorem 17 *When the incompressible condition is satisfied, that is $\nabla \cdot V = 0$, then the flow $\phi_t(x)$ is **measure preserving**, that is:*

$$|\det J_{\phi_t}(x)| = 1, \text{ i.e. } p(\phi_t(x), t) = p_0(x) \quad (45)$$

Definition 3 (Flux) *We find that $V(x, t)p(x, t)$ is actually the flux $\mathcal{F}(x, t)$ of the particle.*

Then the continuity equation can be rewritten as:

$$\frac{\partial p(x, t)}{\partial t} = -\nabla \cdot (\mathcal{F}(x, t)) \quad (46)$$

Then we find that if the flux s.t. $\mathcal{F} = -\frac{1}{2}\nabla \cdot (GG^T p(x, t))$, then $p(x, t)$ describes the diffusion process. This is the famous Fick's Law.

Theorem 18 (Fick's Law) *Fick's Law describes the relationship between the flux $\mathcal{F}(x, t)$ of the particle and the concentration/density $p(x, t)$.*

$$\mathcal{F}(x, t) = -\frac{1}{2}\nabla \cdot (GG^T p(x, t)) \quad (47)$$

Specifically, when $G(X_t, t) = G(t)$ and $Q = I$, we have:

$$\mathcal{F}(x, t) = -\frac{GG^T}{2}\nabla p(x, t) \quad (48)$$

Then

$$\frac{\partial p(x, t)}{\partial t} = \nabla \cdot \left(\frac{GG^T}{2}\nabla p(x, t) \right) \quad (49)$$

4.3 Solution

Note $-\frac{GG^T}{2}$ is actually the diffusion coefficient \mathcal{D} . Then we have the diffusion equation:

$$\frac{\partial p(x, t)}{\partial t} = \nabla \cdot (\mathcal{D}\nabla p(x, t)) \quad (50)$$

with initial condition $p(x, 0) = p_0(x)$. We can use the Fourier Transform to solve this equation.

Theorem 19 *The solution to the diffusion equation is:*

$$\begin{aligned} p(x, t) &= \mathcal{F}^{-1} [\tilde{p}_0(\lambda) \exp(-\lambda^T \mathcal{D} \lambda t)] = (p_0 \star \mathcal{G}_{2t\mathcal{D}})(x) \\ &= \frac{1}{\sqrt{(4\pi t)^d \det(\mathcal{D})}} \int_{\mathcal{R}^d} \left(p_0(\xi) \exp\left(-\frac{1}{4t} (x - \xi)^T \mathcal{D}^{-1} (x - \xi)\right) \right) d\xi \end{aligned} \quad (51)$$

where $\tilde{p}_0(\lambda) = \mathcal{F}(p_0(x))$ is the Fourier Transform of $p_0(x)$. $\mathcal{G}_{2t\mathcal{D}}$ is the Gaussian Kernel with variance $2t\mathcal{D}$.

Proof 9 First, assume the Fourier Transform of $p(x, t)$ is $\tilde{p}(x, t)$:

$$\begin{cases} \tilde{p}(x, t) = \mathcal{F}[p(x, t)] = \int_{\mathcal{R}^d} p(x, t) e^{-i\lambda \cdot x} dx \\ p(x, t) = \mathcal{F}^{-1}[\tilde{p}(x, t)] = \frac{1}{(2\pi)^d} \int_{\mathcal{R}^d} \tilde{p}(x, t) e^{i\lambda \cdot x} dx \end{cases} \quad (52)$$

Then, we have:

$$\begin{cases} \mathcal{F}[\nabla \cdot \mathbf{v}] = i\lambda \cdot \mathcal{F}[\mathbf{v}] \\ \mathcal{F}[\mathcal{D}\nabla p] = i\mathcal{D}\lambda \mathcal{F}[p] \end{cases} \quad (53)$$

Then,

$$\begin{aligned} \mathcal{F}\left[\frac{\partial p}{\partial t}\right] &= \frac{d}{dt} \mathcal{F}[p] = \mathcal{F}[\nabla \cdot (\mathcal{D}\nabla p)] \\ &= i\lambda \cdot \mathcal{F}[\mathcal{D}\nabla p] = -\lambda^T \mathcal{D} \lambda \mathcal{F}[p] \end{aligned} \quad (54)$$

where $\lambda = (\lambda_1, \lambda_2, \dots, \lambda_d)^T$.

Therefore, $\mathcal{F}[p] = \tilde{p}_0 \exp(-\lambda^T \mathcal{D} \lambda t)$. Since $\mathcal{F}[N(x|0, 2t\mathcal{D})] = \exp(-\lambda^T \mathcal{D} \lambda t)$, which gives the theorem.

Remark 4 Specially, 1. When the initial density $p_0(x)$ is $\delta(x - x_0)$, the solution is:

$$p(x, t) = \frac{1}{\sqrt{(4\pi t)^d \det \mathcal{D}}} \exp\left(-\frac{(x - x_0)^T \mathcal{D}^{-1} (x - x_0)}{4t}\right) \sim N(x_0, 2t\mathcal{D}) \quad (55)$$

2. When the initial density $p_0(x)$ is a Gaussian distribution $N(\mu, \Sigma)$, the solution is:

$$p(x, t) = \frac{1}{\sqrt{(2\pi)^d \det(\Sigma + 2t\mathcal{D})}} \exp\left(-\frac{1}{2} (x - \mu)^T (\Sigma + 2t\mathcal{D})^{-1} (x - \mu)\right) \sim N(\mu, \Sigma + 2t\mathcal{D}) \quad (56)$$

(The Fourier transform of (μ, Σ) is $\exp(-i\lambda^T \mu + \frac{1}{2} \lambda^T \Sigma \lambda)$.)

Till here, we can see the insight of diffusion. It is actually a process of smoothing the initial density by the Gaussian Kernel.

5 Variational Auto-Encoder

5.1 Structure

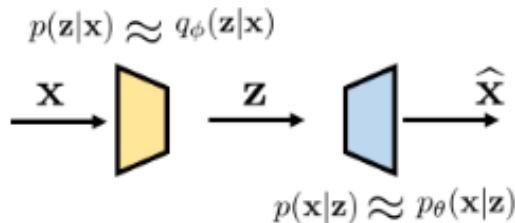


Figure 1: VAE block

where $q_\phi(\mathbf{z} | \mathbf{x})$ is the proxy for $p(\mathbf{z} | \mathbf{x})$, which is also the distribution associated with the encoder. And $p_\theta(\mathbf{x} | \mathbf{z})$ is the proxy for $p(\mathbf{x} | \mathbf{z})$, which is also the distribution associated with the decoder. Like the encoder, the decoder can be parameterized by a deep neural network.

If we treat ϕ and θ as optimization variables, then we need an objective function (or the loss function) so that we can optimize ϕ and θ through training samples.

5.2 Evidence Lower Bound

Definition 4 (Evidence Lower Bound) *The Evidence Lower Bound (ELBO) is defined as:*

$$\text{ELBO}(x) = \mathbf{E}_{q_\phi(z|x)} \left[\log \frac{p(x, z)}{q_\phi(z|x)} \right] \quad (57)$$

Remark 5 *The ELBO is a lower bound of the log-likelihood of the data. It is used to estimate $\log p(x)$.*

$$\begin{aligned} \log p(x) &= \log \int p(x, z) dz \\ &= \log \int \frac{p(x, z)}{q_\phi(z|x)} \cdot q_\phi(z|x) dz \\ &\geq \mathbf{E}_{q_\phi(z|x)} \left[\log \frac{p(x, z)}{q_\phi(z|x)} \right] = \text{ELBO}(x) \end{aligned} \quad (58)$$

Theorem 20 (Decomposition of Log-likelihood) *We have*

$$\log p(x) = \text{ELBO}(x) + \text{KL}(q_\phi(z|x) || p(z|x)) \quad (59)$$

then we can minimize the gap between $\log p(x)$ and ELBO, and the equality hold if and only if $q_\phi(z|x) = p(z|x)$. Since $p(z|x)$ is a delta function,

Proof 10

$$\begin{aligned} \log p(x) &= \log p(x) \int q_\phi(z|x) dz \\ &= \mathbf{E}_{q_\phi(z|x)} [\log p(x)] \\ &= \mathbf{E}_{q_\phi(z|x)} \left[\log \left(\frac{p(x, z)}{p(z|x)} \frac{q_\phi(z|x)}{q_\phi(z|x)} \right) \right] \\ &= \text{ELBO}(x) + \text{KL}(q_\phi(z|x) || p(z|x)) \end{aligned} \quad (60)$$

Theorem 21 *Also, we can rewrite the ELBO as:*

$$\begin{aligned} \text{ELBO}(x) &= \mathbf{E}_{q_\phi(z|x)} [\log p(x|z) + \log p(z) - \log q_\phi(z|x)] \\ &= \mathbf{E}_{q_\phi(z|x)} [\log p(x|z)] - \text{KL}(q_\phi(z|x) || p(z)) \\ &= \mathbf{E}_{q_\phi(z|x)} [\log p_\theta(x|z)] - \text{KL}(q_\phi(z|x) || p(z)) \end{aligned} \quad (61)$$

where the first term determines how good the decoder is, maximizing the likelihood of observing the image, and the latter describes how good the encoder is, minimizing the distance between two distributions.

Definition 5 (The objective of VAE) *The optimization objective of VAE is to maximize the ELBO:*

$$(\phi, \theta) = \operatorname{argmax}_{\phi, \theta} \sum_{x \in X} \text{ELBO}(x) \quad (62)$$

where X is the training set.

6 Diffusion Model

6.1 DDPM

DDPM is like splitting the encoder and decoder of VAE into controllable parts. For each training data point $x_0 \sim p_{data}$, then a discrete Markov chain $\{x_1, \dots, x_N\}$ is constructed by transition function:

$$p(x_i | x_{i-1}) = \mathcal{N}(x_i | \sqrt{1 - \beta_i} x_{i-1}, \beta_i I) \quad (63)$$

Then we can get

$$p_{\alpha_i}(x_i | x_0) = \mathcal{N}(x_i | \sqrt{\alpha_i} x_0, (1 - \alpha_i) I), \alpha_i = \prod_{j=1}^i (1 - \beta_j) \quad (64)$$

Hence we need to train the ELBO:

$$\text{ELBO}(x) = \sum_{i=1}^N (1 - \alpha_i) \mathbf{E}_{p_{data}(x)} \left[\mathbf{E}_{p_{\alpha_i}(\hat{x}|x)} [||s_\theta(\hat{x}, i) - \nabla_{\hat{x}} \log p_{\alpha_i}(\hat{x}|x)||] \right] \quad (65)$$

Then do the reverse Markov chain.

This is clearly a discrete version. Then we consider the continuous version. We consider linear SDE having the form:

$$dX_t = (a(t)X_t + b(t))dt + g(t)dW_t \quad (66)$$

where $X_t \in \mathcal{R}^d, W_t \in \mathcal{R}^m$ with diffusion factor $Q \in \mathcal{R}^{m \times m}$, then $a(t) \in \mathcal{R}^{d \times d}, b(t) \in \mathcal{R}^d, g(t) \in \mathcal{R}^{d \times m}$. By Euler Maruyama method, it can be approximated By

$$\begin{aligned} X_{t+s} &= X_t + (a(t)X_t + b(t))s + g(t)\sqrt{sQ}\xi \\ &= (1 + a(t)s)X_t + b(t)s + g(t)\sqrt{sQ}\xi \end{aligned} \quad (67)$$

where $\xi \sim N(0, I_m)$. Usually we need to consider the expectation, variance and distribution of X_t . But the stochastic value of X_t is dependent of x_0 . Then first we consider

$$\begin{aligned} E[X_{t+s}|X_0] - E[X_t|X_0] &\approx (a(t)E[X_t|X_0] + b(t))s + g(t)\sqrt{sQ}E[\xi] \\ &= (a(t)E[X_t|X_0] + b(t))s. \end{aligned} \quad (68)$$

Note $e(t) = E[X_t|X_0]$, then

$$e'(t) = \lim_{s \rightarrow 0} \frac{E[X_{t+s}|X_0] - E[X_t|X_0]}{s} = a(t) \cdot e(t) + b(t). \quad e(0) = X_0. \quad (69)$$

which is an ODE system, having solution

$$e(t) = e^{\int_0^t a(s)ds} \cdot \left(X_0 + \int_0^t e^{-\int_0^s a(r)dr} b(s)ds \right) \quad (70)$$

Therefore

$$\begin{aligned} E[X_t] &= E[E[X_t|X_0]] = E[e(t)] \\ &= e^{\int_0^t a(s)ds} \cdot \left(E[X_0] + \int_0^t e^{-\int_0^s a(r)dr} b(s)ds \right) \end{aligned} \quad (71)$$

Similarly, Note $\text{Var}(X_0|X_0) = v(t)$: then $\text{Var}(X_{t+s}|X_0) = (1 + sa(t))^2 \text{Var}(X_t|X_0) + sgQg^\top$. Then

$$\begin{aligned} V'(t) &= \lim_{s \rightarrow 0} \frac{\text{Var}(X_{t+s}|X_0) - \text{Var}(X_t|X_0)}{s} \\ &= [(a^2(t)s + 2a(t))v(t) + g^2(t)]|_{s \rightarrow 0} \\ &= 2\alpha(t)V(t) + g(t)Qg^\top(t), \quad V(0) = 0 \end{aligned} \quad (72)$$

Solution is:

$$v(t) = e^{\int_0^t 2a(s)ds} \cdot \left(\int_0^t e^{-\int_0^s 2a(r)dr} g(s)Qg^\top(s)ds \right) \quad (73)$$

By law of total variance:

$$\begin{aligned} \text{Var}(X_t) &= E[X_t^2] - E^2[X_t] = E[E[X_t^2|X_0]] - E^2[X_t] \\ &= E[\text{Var}(X_t|X_0) + E^2[X_t|X_0]] - E^2[X_t] \\ &= E[\text{Var}(X_t|X_0)] + E[E^2[X_t|X_0]] - E^2[E[X_t|X_0]] \\ &= E[\text{Var}(X_t|X_0)] + \text{Var}(E[X_t|X_0]) \end{aligned} \quad (74)$$

then

$$\begin{aligned} \text{Var}(X_t) &= E[V(t)] + \text{Var}(e(t)) \\ &= e^{\int_0^t 2a(s)ds} \cdot \left(\int_0^t e^{-\int_0^s 2a(r)dr} g(s)Qg^\top(s)ds \right) + e^{\int_0^t 2a(s)ds} \cdot \text{Var}(X_0). \end{aligned} \quad (75)$$

We have the following theorem which is crucial for diffusion models. Usually, we assume $Q = I_m$.

Theorem 22 If $X_{t+s} = (1 + a(t)s)X_t + b(t)s + g(t)\sqrt{sQ}\xi$ then $X_t|X_0 \sim N(E[X_t|X_0], \text{Var}(X_t|X_0))$, where $E[X_t|X_0] = e(t)$, $\text{Var}(X_t|X_0) = V(t)$.

It should be noted that $e(t)$ is related to X_0 and t , while $V(t)$ only depends on t !

Next, we will see how the above formula can be applied to diffusion models. There are three frameworks to build SDEs for diffusion models, VP, VE and sub-VP.

Definition 6 Noise function $\beta(t)$. s.t. $\beta(0) = 0; \beta'(t) \geq 0; \beta(t) \rightarrow \infty$ as $t \rightarrow \infty$

Variance Preserving (VP) SDE

So if we have diffusion model like:

$$\begin{aligned} X_{t_{i+1}} &= \sqrt{1 - (\beta(t_{i+1}) - \beta(t_i))} X_{t_i} + \sqrt{(\beta(t_{i+1}) - \beta(t_i))} \xi \\ &= \sqrt{1 - \Delta\beta(t_i)} X_{t_i} + \sqrt{\Delta\beta(t_i)} \xi \end{aligned} \quad (76)$$

Then the conditional distribution is given by:

$$q(X_{t_{i+1}}|X_{t_i}) = N(x_{t_{i+1}}; \sqrt{1 - \Delta\beta(t_i)} X_{t_i}, \Delta\beta(t_i)) \quad (77)$$

Then we need to estimate θ drift term f and diffusion term g :

$$\begin{aligned} f(x, t) &= \lim_{h \rightarrow 0} \frac{E[X_{t+h} - X_t | X_t = x]}{h} \\ &= \lim_{h \rightarrow 0} \frac{x\sqrt{1 - \Delta\beta(t)} - x}{h} = -\frac{x}{2}\beta'(t). \\ g(t) &= \sqrt{\lim_{h \rightarrow 0} \frac{N[X_{t+h}|X_t = x]}{h}} = \sqrt{\lim_{h \rightarrow 0} \frac{\beta(t+h) - \beta(t)}{h}} = \sqrt{\beta'(t)} \end{aligned} \quad (78)$$

Then the model can be written as $dx = -\frac{x}{2}\beta'(t)dt + \sqrt{\beta'(t)}dW_t$

Then by Theorem 22 we have

$$\begin{cases} E[X_t|X_0] = X_0 e^{\int_0^t -\frac{1}{2}\beta'(s)ds} = X_0 e^{-\frac{1}{2}\beta(t)} \\ E[X_t] = E[X_0] e^{-\frac{1}{2}\beta(t)} \\ V(X_t|X_0) = \int_0^t e^{\int_0^s \beta'(r)dr} \beta'(s) ds \cdot e^{-\beta(t)} = 1 - e^{-\beta(t)} \\ V(X_t) = 1 - e^{-\beta(t)} + V(X_0) e^{-\beta(t)} = 1 + (V(X_0) - 1) e^{-\beta(t)}. \end{cases} \quad (79)$$

So as $t \rightarrow \infty, \beta(t) \rightarrow \infty$, then $E \rightarrow 0, V \rightarrow 1$, i.e. $X_t|X_0 \sim N(E[X_t|X_0], \text{Var}[X_t|X_0]) \rightarrow N(0, 1)$ as $t \rightarrow \infty$.

Variance-Exploding SDE

Here is the model: $X_{t+h} = X_t + \sqrt{\Delta\beta(t)}\xi$

Similarly we can compute the $f(x, t) \equiv 0$ and $g(t) = \sqrt{\beta'(t)}$. Hence

$$\begin{cases} E[X_0|X_0] = X_0 \\ E[X_t] = E[X_0] \\ V(X_t|X_0) = \int_0^t e^{\int_0^s 0dr} \beta'(s) ds = \beta(t) \\ V(X_t) = V[X_0] + \beta(t) \end{cases} \quad (80)$$

So the expectation value is constant and the variance is increasing monotonical.

If we rescale X_t as $Y_t = \frac{X_t}{\sqrt{\beta(t)}}$, then $Y_t \rightarrow N(0, 1), t \rightarrow \infty$.

Sub-VP SPE

Here, we set the drift and diffusion term as

$$\begin{aligned} f(x, t) &= -\frac{1}{2}\beta'(t) \\ g(t) &= \sqrt{\beta'(t) (1 - e^{-2\beta(t)})} \end{aligned} \quad (81)$$

As the same, we can compute that.

$$\begin{cases} E[X_t|X_0] = X_0 e^{-\frac{1}{2}\beta(t)} \\ E[X_t] = E[X_0] e^{-\frac{1}{2}\beta(t)} \\ V(X_t|X_0) = (1 - e^{-\beta(t)})^2 \\ V(X_t) = (1 - e^{-\beta(t)})^2 + V(X_0) e^{-\beta(t)}. \end{cases} \quad (82)$$

We can find out that the variance is always smaller than that of VP SDE.

Remark 6 To sum up, finally we hope that X_t converges to a normal distribution by choosing different drift and diffusion functions. For generative model, the goal is to sample from a Data distribution p_{data} . We have known that if we set the initial distribution $p_0(x_0) = p(X_0 = x_0) \sim p_{data}$, then after $t = T$, the distribution of X_t is tend to be $N(0, 1)$ under certain conditions.

So the idea is backward: if we sample from $X_T \sim N(0, 1)$, and then run SDE backwards, could we get the initial distribution?

Assume we have forward SDE: from $X_0 \sim p_0, X_T \sim p_T$,

$$dX_t = f(X_t, t)dt + G(t)dW_t \quad (83)$$

Then we define the reverse SDE as: from $X_T \sim p_T$,

$$d\bar{X}_t = \bar{f}(\bar{X}_t, t)dt + \bar{G}(t)d\bar{W}_t \quad (84)$$

where \bar{W}_t is Brownian Motion runs backward in time, i.e. $\bar{W}_{t-s} - \bar{W}_t$ is independent of \bar{W}_t . We can approximate by EM:

$$\bar{X}_{t-s} - \bar{X}_t = -s\bar{f}(\bar{X}_t, t) + \sqrt{s}\bar{G}(t)\xi \quad (85)$$

So the problem is: If given f, G , are there \bar{f}, \bar{G} s.t. the reverse time diffusion process \bar{X}_t has the same distribution as the forward process X_t ? Yes!

Theorem 23 The reverse SDE with \bar{f}, \bar{G} having the following form has the same distribution as the forward SDE 83:

$$\begin{cases} \bar{f}(x, t) = f(x, t) - GG^T \nabla_x \log p_t(x) \\ \bar{G} = G(t) \end{cases} \quad (86)$$

i.e.

$$d\bar{X}_t = [f(\bar{X}_t, t) - GG^T \nabla_x \log p_t(x_t)] dt + G(t)d\bar{W}_t \quad (87)$$

Proof 11 The proof is skipped.

This theroem allows us to learn how to generate samples from p_{data} .

Algorithm 1 :

Step1. Select $f(x, t)$ and $g(t)$ with affine drift coefficients s.t. $X_T \sim N(0, 1)$

Step2. Train a network $s_\theta(x, t) = \frac{\partial}{\partial x} \log p_t(x)$ where $p_t(x) = p(X_t = x)$ is the forward distribution.

Step3. Sample X_T from $N(0, 1)$, then run reverse SDE from T to 0:

$$\bar{X}_{t-s} = \bar{X}_t + s [g^2(t)s_\theta(\bar{X}_t, t) - f(\bar{X}_t, t)] + \sqrt{s}g(t)\xi \quad (88)$$

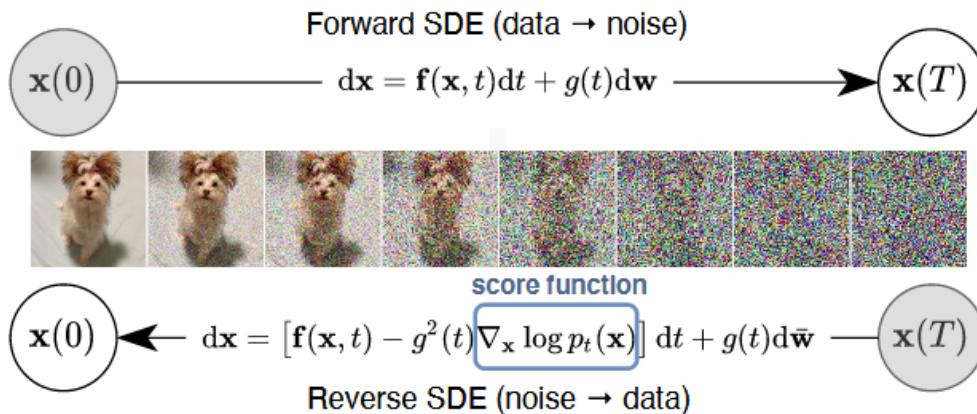


Figure 2: score-based generative model

6.2 Score Matching

The most difficult question on how to obtain $\nabla_x \log p(x)$ because it solves FPK equation.

6.2.1 Explicit Score Matching

Suppose we have a set of samples x_1, x_2, \dots, x_n from the data distribution $p_{data}(x)$. A classical way is to consider the kernel density estimation $q(x)$ of $p(x)$:

$$q(x) = \frac{1}{n} \sum_{i=1}^n K(x - x_i) \quad (89)$$

where $K(x)$ is the kernel function. Since $q(x)$ is an approximation to p_{data} . We can define a loss function to train a network:

$$\begin{aligned} \mathcal{L}_\theta &= \mathbf{E}_{x \sim p(x)} \left[\|s_\theta(x) - \nabla_x \log p(x)\|^2 \right] \\ &\approx \mathbf{E}_{x \sim q(x)} \left[\|s_\theta(x) - \nabla_x \log q(x)\|^2 \right] \\ &= \int \|s_\theta(x) - \nabla_x \log q(x)\|^2 q(x) dx \\ &\approx \frac{1}{n} \sum_{i=1}^n \int \|s_\theta(x) - \nabla_x \log q(x)\|^2 K(x - x_i) dx \end{aligned} \quad (90)$$

However, when the number of samples is limited, the estimation $\nabla_x \log q(x)$ is not accurate.

6.2.2 Implicit Score Matching

6.2.3 Denoising Score Matching

Normally we can define the loss function as follows:

$$\begin{aligned} L_\theta &= \frac{1}{T} \int_0^T \lambda(t) \mathbf{E}_{x_0 \sim p_{data}} \left[\mathbf{E}_{x_t \sim p_{t|0}(x_t|x_0)} \left[\|s_\theta(x_t, t) - \nabla_{x_t} \log p_t(x_t)\|^2 \right] \right] dt \\ &= \mathbf{E}_{t \sim U(0, T)} \left[\lambda(t) \mathbf{E}_{x_0 \sim p_{data}} \left[\mathbf{E}_{x_t \sim p_{t|0}(x_t|x_0)} \left[\|s_\theta(x_t, t) - \nabla_{x_t} \log p_t(x_t)\|^2 \right] \right] \right] \end{aligned} \quad (91)$$

It should be clarified that $p_{t|0}(x_t|x_0) = p(X_t = x_t | X_0 = x_0)$. So

$$p_t(x_t) = \int p_{t|0}(x_t|x_0) p_0(x_0) dx_0 = \mathbf{E}_{x_0 \sim p_{data}} [p_{t|0}(x_t|x_0)]$$

where $p_t(x) = p(X_t = x)$, $p_{t|0}(x|y) = p(X_t = x | X_0 = y)$. Then

$$\begin{aligned} \nabla \log p_t(x) &= \frac{1}{p_t(x)} \nabla p_t(x). \\ &= \frac{1}{p_t(x)} \nabla \int p_{t|0}(x|y) p_0(y) dy \\ &= \frac{1}{p_t(x)} \int \nabla p_{t|0}(x|y) p_0(y) dy \\ &= \frac{1}{p_t(x)} \int \frac{\nabla p_{t|0}(x|y)}{p_{t|0}(x|y)} p_0(y) \cdot p_{t|0}(x|y) dy \\ &= \int \nabla_x \log(p_{t|0}(x|y)) \cdot p_{0|t}(y|x) dy \\ &= \mathbf{E}_{y \sim p_{0|t}(y|x)} [\nabla_x \log(p_{t|0}(x|y))] \end{aligned} \quad (92)$$

Where we have used the following lemma:

Lemma 1

$$\mathbf{E}_{x_0 \sim p_0} \left[\mathbf{E}_{x_t \sim p_{t|0}(\cdot|x_0)} \left[\mathbf{E}_{x'_0 \sim p_{0|t}(\cdot|x_t)} [f(x_t, x'_0)] \right] \right] = \mathbf{E}_{x_0 \sim p_0} \left[\mathbf{E}_{x_t \sim p_{t|0}(\cdot|x_0)} [f(x_t, x_0)] \right] \quad (93)$$

Proof 12 *Easy to prove.*

Then we can rewrite the loss function as:

$$\begin{aligned}
L_\theta &= \mathbb{E}_{t \sim U(0,T)} \left[\lambda(t) \mathbb{E}_{x_0 \sim p_{data}} \left[\mathbb{E}_{x_t \sim p_{t|0}(x_t|x_0)} \left[\|S_\theta(x_t, t) - \nabla_{x_t} \log p_t(x_t)\|^2 \right] \right] \right] \\
&\leq \mathbb{E}_{t \sim U(0,T)} \left[\lambda(t) \mathbb{E}_{x_0 \sim p_{data}} \left[\mathbb{E}_{x_t \sim p_{t|0}(x_t|x_0)} \left[\mathbb{E}_{y \sim p_{data}} \left[\|S_\theta(x_t, t) - \nabla_{x_t} \log(p_{t|0}(x_t|y))\|^2 \right] \right] \right] \right] \\
&= \mathbb{E}_{t \sim U(0,T)} \left[\lambda(t) \mathbb{E}_{x_0 \sim p_{data}} \left[\mathbb{E}_{x_t \sim p_{t|0}(x_t|x_0)} \left[\|S_\theta(x_t, t) - \nabla_{x_t} \log(p_{t|0}(x_t|x_0))\|^2 \right] \right] \right]
\end{aligned} \tag{94}$$

Since $p_{t|0}(x_t|x_0) = p(X_t = x_t|X_0 = x_0)$ has been discussed:

$$p_{t|0}(x_t|x_0) \sim N(x_t; E[X_t = x_t|X_0 = x_0], \text{Var}(X_t = x_t|X_0 = x_0)).$$

Then by theorem 22, x can be written as $x = e(t, X_0) + \sqrt{V(t)}\xi$, where $\xi \sim N(0, 1)$, then the score function is:

$$\frac{\partial}{\partial x} \log p_{t|0}(x|x_0) = -\frac{x - E_{t|0}[x|x_0]}{\text{Var}_{t|0}(x|x_0)} = -\frac{x - e(t, X_0)}{V(t)} \sim -N\left(0, \frac{1}{V(t)}\right) \tag{95}$$

So

$$\begin{aligned}
L_\theta &= \mathbb{E}_{t \sim U(0,T)} \left[\lambda(t) \mathbb{E}_{x_0 \sim p_{data}} \left[\mathbb{E}_{\xi \sim N(0,1)} \left[\left\| s_\theta\left(\sqrt{V(t)}\xi + e(t, X_0), t\right) + \frac{\xi}{\sqrt{V(t)}} \right\|^2 \right] \right] \right] \\
&= \mathbb{E}_{t \sim U(0,T)} \left[\lambda(t) \mathbb{E}_{x_0 \sim p_{data}} \left[\frac{1}{V(t)} \mathbb{E}_{\xi \sim N(0,1)} \left[\left\| \xi_\theta\left(\sqrt{V(t)}\xi + e(t, X_0), t\right) - \xi \right\|^2 \right] \right] \right]
\end{aligned} \tag{96}$$

where $\xi_\theta = -\sqrt{V(t)}s_\theta$ is called denoising network.

6.3 Denoising Diffusion

6.3.1 With Classifier Guidance

Though we can produce pictures by sampling from normal distribution, we still cannot control what we will generate. What we want to do is something like: "Give me the pictures of number 6", then the model can sample from the normal distribution and do the denoising to generate pics of 6.

Usually, we can do something like: train a model for every class label. This do make the model smaller, but increases number of models. Think about it, when the label is TEXT, it is impossible to train a model for each sentences.

So, the initial distribution is $p_0(x|y)$ given the label y . Similarly, we will convert the data distribution $p_{data}(x|y)$ to final distribution, normal distribution expected. Then we SDE becomes: $X_t \sim p_t(x|y)$

$$\begin{aligned}
p_t(x|y) &= p(X_t = x|y) = \frac{p(y|X_t = x)p(X_t = x)}{p(y)} \\
&\Rightarrow \log(p_t(x|y)) = \log(p(y|X_t = x)) + \log(p(X_t = x)) - \log(p(y)) \\
&\Rightarrow \nabla_x \log(p_t(x|y)) = \nabla_x \log(p(y|X_t = x)) + \nabla_x \log(p(X_t = x))
\end{aligned} \tag{97}$$

We have finished training $\nabla_x \log(p(X_t = x))$ in sampling. Then we need to estimate $\nabla_x \log(p(y|X_t = x))$. This is the conditional protability, we end up with a sharp factor s: $p'(y|X_t = x)$, then:

$$\nabla_x \log(p_t(x|y)) = S \nabla_x \log(p(y|x_t = x)) + \nabla_x \log(p(x_t = x)) \tag{98}$$

Note $\omega_\theta(y|x, t)$ to learn $S \nabla_x \log(p(y|X_t = x))$

6.3.2 Classifier Guidance Free

$$\begin{aligned}
&\gamma \nabla_x \log(p(y|X_t = x)) \\
&= \gamma (\nabla_x \log(p(X_t = x|y)) - \nabla_x \log(p_t(x)))
\end{aligned} \tag{99}$$

Then

$$\begin{aligned}
&\nabla_x \log_\gamma(p_t(x|y)) \\
&= (1 - \gamma) \nabla_x \log(p_t(x)) + \gamma \nabla_x \log(p(X_t = x|y))
\end{aligned} \tag{100}$$

Hence we only need one conditional denoising network, and using null condition to represent the unconditional model.

7 Flow Matching

7.1 FPK Equation

We have discussed the FPK Equation in 'learnsde'.

Theorem 24 (Fokken-Planck-Kolmogorov equation) *The density function $p(x, t)$ of X_t s.t.*

$$dX_t = f(X_t, t)dt + G(X_t, t)dW_t \quad (101)$$

solves the PDE:

$$\frac{\partial p(x, t)}{\partial t} = - \sum_i \frac{\partial}{\partial x_i} [f_i(x, t)p(x, t)] + \frac{1}{2} \sum_{i,j} \frac{\partial^2}{\partial x_i \partial x_j} [(GQG^\top)_{ij} p(x, t)] \quad (102)$$

The PDE is called FPK equation / forward Kolmogorov equation.

It can be rewritten as:

$$\begin{aligned} \frac{\partial p(x, t)}{\partial t} &= -\nabla \cdot [f(x, t)p(x, t)] + \frac{1}{2} \nabla^2 \cdot [(GQG^\top) p(x, t)] \\ &= -\nabla \cdot \left[f(x, t)p(x, t) - \frac{1}{2} \nabla \cdot [(GQG^\top) p(x, t)] \right] \end{aligned} \quad (103)$$

Here, if we only consider $G(X_t, t) = g(t)$, then we notice that $M = GQG^\top$ is independent of X_t , so we can write:

$$\begin{aligned} \frac{\partial p(x, t)}{\partial t} &= -\nabla \cdot \left(fp - \frac{1}{2} \nabla \cdot (Mp) \right) \\ &= -\nabla \cdot \left(fp - \frac{1}{2} M \nabla p \right) \\ &= -\nabla \cdot \left[\left(f - \frac{1}{2} M \frac{\nabla p}{p} \right) p \right] \\ &= -\nabla \cdot \left[\left(f - \frac{1}{2} M \nabla \log p \right) p \right] \end{aligned} \quad (104)$$

So we find out that if we have an ODE s.t. $dZ_t = F(Z_t, t)dt$ with $Z_0 \sim p_0$, instead of a SDE, then by FPK equation, the density $p(z, t)$ satisfies:

$$\frac{\partial p(z, t)}{\partial t} = -\nabla \cdot (F(z, t)p(z, t)) \quad (105)$$

So if we set $F(z, t) = f(z, t) - \frac{1}{2} M(t) \nabla \log p(z, t)$, then $p(z, t)$ is exactly like the density $p(x, t)$ of X_t in SDE. So theoretically, we can do the diffusion like reverse ode!

This is the topic discussed in 'Probability Flow'.

7.2 Probability Flow

Define a flow $\phi : \mathcal{R}^d \times [0, 1] \rightarrow \mathcal{R}^d$ is a flow generated by a vector field $v : \mathcal{R}^d \times [0, 1] \rightarrow \mathcal{R}^d$ i.e.

$$\begin{cases} \frac{\partial \phi(x, t)}{\partial t} = v(\phi(x, t), t) \\ \phi(0, x) = x \end{cases} \quad (106)$$

The flow means that under the vector field v , if the initial point is x , then the flow push the point after time t to $\phi(x, t)$. That is ϕ gives the evolution trajectory of x under the vector field v . So normally, we can consider the flow $\phi(x, t)$ as X_t in SDE:

$$dX_t = v(X_t, t)dt + 0dW_t \quad (107)$$

with $X_0 = x$. It turns out that it is actually an ODE, a SDE without diffusion term. Similar to SDE, if $X_0 = x \sim p_0(x)$, we have the probability density $p(x, t)$ satisfies FPK equation:

$$\frac{\partial p(x, t)}{\partial t} = -\nabla \cdot (v(x, t)p(x, t)) \quad (108)$$

which is a special case of FPK equation, called **Continuity Equation**. So, typically, we can solution to an ODE is a flow.

So the objective of Flow Matching Model can be described as: Let $p_t(x)$ be the density with initial $p_0(x)$, which is designed to be a simple distribution, like normal distribution. So let $p_1(x)$ be the approximation equal in distribution p_{data} . Then we need to design a flow to match the flow s.t. p_1 can properly approximate p_{data} .

8 Semilinear Stochastic PDEs

8.1 Semilinear SPDE

Then we come to the time-dependent SPDE. We study the stochastic semilinear evolution equation:

$$du = [\Delta u + f(u)]dt + G(u)dW(t, x) \quad (109)$$

Definition 7 (Semilinear SPDE) *Simmlar to normal time-dependent PDE, we treat SPDE like this as semi-linear SODEs on a Hilbert space, like*

$$du = [-Au + f(u)]dt + G(u)dW(t) \quad (110)$$

where $-A$ is a linear operator that generates a semigroup $S(t) = e^{-tA}$.

Example 2 (Phase-field model)

$$du = [\epsilon \Delta u + u - u^3]dt + \sigma dW(t, x) \quad (111)$$

Example 3 (Fluid Flow)

$$\begin{aligned} u_t &= \epsilon \Delta u - \nabla p - (u \cdot \nabla)u \\ \nabla \cdot u &= 0 \end{aligned} \quad (112)$$

So like we deal with integration of stochastic process like Itos or stratonovich, we need to generalize the Brownian Motion by introducing spatial variable to $W(t)$. Here we define Q-Wiener Process.

8.2 Q wiener process

First, we assume U is a Hilbert space. And $(\Omega, \mathbf{F}, \mathbf{F}_t, \mathbb{R})$ is a filtered probability space.

Definition 8 (Q) $Q \in \mathcal{L}(U)$ is non-negative definite and symmetric. Further, Q has an orthonormal basis $\{\mathcal{X}_j : j \in \mathcal{N}\}$ of eigenfunctions with corresponding eigenvalues $q_j \geq 0$ such that $\sum_{j \in \mathcal{N}} q_j < \infty$ (i.e., Q is of trace class).

Definition 9 (Q-Wiener Process) A U -valued stochastic process $\{W(t) : t \geq 0\}$ is Q -Wiener process if

- $W(0) = 0$ a.s.
- $W(t)$ is a continuous function $\mathbb{R}^+ \rightarrow U$, for each $\omega \in \Omega$.
- $W(t)$ is \mathcal{F}_t -adapted and $W(t) - W(s)$ is independent of \mathcal{F}_s for $s \leq t$
- $W(t) - W(s) \sim N(0, (t-s)Q)$ for all $0 \leq s \leq t$

Theorem 25 (Q-Wiener Process) Assume we have Q defined in 8. Then, $W(t)$ is a Q -Wiener process if and only if

$$W(t) = \sum_{j=1}^{\infty} \sqrt{q_j} \mathcal{X}_j \beta_j(t) \quad (113)$$

which is converges in $L^2(\Omega, C([0, T], U))$ and $\beta_j(t)$ are iid \mathcal{F}_t -Brownian motions and the series converges in $L^2(\Omega, U)$.

Theorem 26 ($H_{\text{per}}^r(0, a)$ -valued process) ...

Theorem 27 ($H_0^r(0, a)$ -valued process) ...

So, in place of $L^2(D)$, we develop the theory on a separable Hilbert space U with norm $\|\cdot\|_U$ and inner product $\langle \cdot, \cdot \rangle_U$ and define the Q-Wiener process $W(t) : t \geq 0$ as a U -valued process.

8.3 Cylindrical Wiener Process

We mention the important case of $Q = I$, which is not trace class on an infinite-dimensional space U (as $q_j = 1$ for all j) so that the series does not converge in $L^2(\Omega, U)$. To extend the definition of a Q -Wiener process, we introduce the cylindrical Wiener process.

The key point is to introduce a second space U_1 such that $U \subset U_1$ and $Q = I$ is a trace class operator when extended to U_1 .

Then we can define cylindrical Wiener process:

Definition 10 (Cylindrical Wiener Process) *Let U be a separable Hilbert space. The cylindrical Wiener process (also called space-time white noise) is the U -valued stochastic process $W(t)$ defined by*

$$W(t) = \sum_{j=1}^{\infty} \mathcal{X}_j \beta_j(t)$$

where $\{\mathcal{X}_j\}$ is any orthonormal basis of U and $\beta_j(t)$ are iid \mathcal{F}_t -Brownian motions.

Theorem 28 *If for the second Hilbert space U_1 , and the inclusion map $\mathcal{I} : U \rightarrow U_1$ is Hilbert-Schmidt. Then, the cylindrical Wiener process is a Q -Wiener process well-defined on U_1 (Converges in $L^2(U, U_1)$).*

8.4 Ito integral solution

Here we consider the Ito integral $\int_0^t B(s) dW(s)$ for a Q -Wiener process $W(s)$. Since dW_t takes value in Hilbert space U , and we treat SPDE in Hilbert space H , the integral will also take value in Hilbert space H .

Hence, $B(s)$ should be $\mathcal{L}_0^2(U_0, H)$ -valued process, where $U_0 \subset U$ known as Cameron-Martin space. So, $B(s)$ is an operator from U_0 to H . Then, we consider the set of operator B .

Definition 11 (L_0^2 space) *Let $U_0 := \{Q^{\frac{1}{2}}u : u \in U\}$, the set of linear operators $B : U_0 \rightarrow H$ is noted as L_0^2 s.t.*

$$\|B\|_{L_0^2} := \left(\sum_{j=1}^{\infty} \|BQ^{\frac{1}{2}}\mathcal{X}_j\|^2 \right)^{\frac{1}{2}} = \|BQ^{\frac{1}{2}}\|_{\text{HS}(U_0, H)} < \infty \quad (114)$$

Remark 7 *If G is invertible, L_0^2 is the space of Hilbert-Schmidt operators $\text{HS}(U_0, H)$.*

Definition 12 *The stochastic integral can be defined by*

$$\int_0^t B(s) dW(s) := \sum_{j=1}^{\infty} \int_0^t B(s) \sqrt{q_j} \mathcal{X}_j d\beta_j(s) \quad (115)$$

So, we can have the truncated form:

$$\int_0^t B(s) dW^J(s) = \sum_{j=1}^J \int_0^t B(s) \sqrt{q_j} \mathcal{X}_j d\beta_j(s) \quad (116)$$

8.5 Semilinear SPDE

Consider the semilinear SPDE:

$$du = [-Au + f(u)]dt + G(u)dW(t) \quad (117)$$

given the initial condition $u_0 \in H$ and $A : \mathcal{D} \subset H \rightarrow H$ is a linear operator, $f : H \rightarrow H$ and $G : H \rightarrow L_0^2$.

Example 4 *Consider the stochastic heat equation:*

$$du = \Delta u dt + \sigma dW(t, x), u(0, x) = u_0(x) \in L^2(D) \quad (118)$$

where D is a bounded domain in \mathbb{R}^d and σ is a constant. Also, homogeneous Dirichlet boundary condition is imposed on D . Hence,

$$H = U = L^2(D), f(u) = 0, G(u) = \sigma I \quad (119)$$

We see that $A = -\Delta$ with domain $\mathcal{D}(A) = H^2(D) \cap H_0^1(D)$.

In the deterministic setting of PDEs, there are a number of different concepts of solution. Here is the same for SPDEs. We can also define strong solution, weak solution and mild solution.

Definition 13 (strong solution) A predictable H -valued process $\{u(t) : t \in [0, T]\}$ is called a strong solution if

$$u(t) = u_0 + \int_0^t [-Au(s) + f(u(s))]ds + \int_0^t G(u(s))dW(s), \quad \forall t \in [0, T] \quad (120)$$

Definition 14 (weak solution) A predictable H -valued process $\{u(t) : t \in [0, T]\}$ is called a weak solution if

$$\langle u(t), v \rangle = \langle u_0, v \rangle + \int_0^t [-\langle u(s), Av \rangle + \langle f(u(s)), v \rangle]ds + \int_0^t \langle G(u(s))dW(s), v \rangle, \quad \forall t \in [0, T], v \in \mathcal{D}(A) \quad (121)$$

where

$$\int_0^t \langle G(u(s))dW(s), v \rangle := \sum_{j=1}^{\infty} \int_0^t \langle G(u(s))\sqrt{q_j}\mathcal{X}_j, v \rangle d\beta_j(s).$$

Definition 15 (mild solution) A predictable H -valued process $\{u(t) : t \in [0, T]\}$ is called a mild solution if for $t \in [0, T]$

$$u(t) = e^{-tA}u_0 + \int_0^t e^{-(t-s)A}f(u(s))ds + \int_0^t e^{-(t-s)A}G(u(s))dW(s),$$

where e^{-tA} is the semigroup generated by $-A$. The right hand side is also called stochastic convolution.

Example 5 (stochastic heat equation in one dimension) Consider the weak solution of 1D heat SPDE with $D = (0, \pi)$, so that $-A$ has eigenfunctions $\phi_j(x) = \sqrt{2/\pi} \sin(jx)$ and eigenvalues $\lambda_j = j^2$ for $j \in \mathbb{N}$. Suppose that $W(t)$ is a Q -Wiener process and the eigenfunctions \mathcal{X}_j of Q are the same as the eigenfunctions ϕ_j of A . A weak solution satisfies: $\forall v \in \mathcal{D}(A)$,

$$\begin{aligned} \langle u(t), v \rangle_{L^2(0, \pi)} &= \langle u_0, v \rangle_{L^2(0, \pi)} + \int_0^t \langle -u(s), Av \rangle_{L^2(0, \pi)} ds \\ &\quad + \sum_{j=1}^{\infty} \int_0^t \sigma \sqrt{q_j} \langle \phi_j, v \rangle_{L^2(0, \pi)} d\beta_j(s) \end{aligned} \quad (122)$$

Assume $u(t) = \sum_{j=1}^{\infty} \hat{u}_j(t) \phi_j$ for $\hat{u}_j(t) := \langle u(t), \phi_j \rangle_{L^2(0, \pi)}$. Take $v = \phi_j$, we have

$$\hat{u}_j(t) = \hat{u}_j(0) + \int_0^t (-\lambda_j) \hat{u}_j(s) ds + \int_0^t \sigma \sqrt{q_j} d\beta_j(s). \quad (123)$$

Hence, $\hat{u}_j(t)$ satisfies the SODE

$$d\hat{u}_j = -\lambda_j \hat{u}_j dt + \sigma \sqrt{q_j} d\beta_j(t) \quad (124)$$

Therefore, each coefficient $\hat{u}_j(t)$ is an Ornstein-Uhlenbeck (OU) process (see Examples 8.1 and 8.21), which is a Gaussian process with variance

$$\text{Var}(\hat{u}_j(t)) = \frac{\sigma^2 q_j}{2\lambda_j} (1 - e^{-2\lambda_j t}) \quad (125)$$

For initial data $u_0 = 0$, we obtain, by the Parseval identity (1.43),

$$\|u(t)\|_{L^2(\Omega, L^2(0, \pi))}^2 = \mathbb{E} \left[\sum_{j=1}^{\infty} |\hat{u}_j(t)|^2 \right] = \sum_{j=1}^{\infty} \frac{\sigma^2 q_j}{2\lambda_j} (1 - e^{-2\lambda_j t}). \quad (126)$$

References

- [1] Stanley H. Chan. Tutorial on diffusion models for imaging and vision, 2024.
- [2] Peter Holderrieth. The fokker-planck equation and diffusion models,, 2023. <https://www.peterholderrieth.com/blog/2023/The-Fokker-Planck-Equation-and-Diffusion-Models/>.
- [3] Yaron Lipman, Marton Havasi, Peter Holderrieth, Neta Shaul, Matt Le, Brian Karrer, Ricky T. Q. Chen, David Lopez-Paz, Heli Ben-Hamu, and Itai Gat. Flow matching guide and code, 2024.
- [4] Gabriel J. Lord, Catherine E. Powell, and Tony Shardlow. *An Introduction to Computational Stochastic PDEs*. Cambridge University Press, Cambridge, 2014.
- [5] Bernt Øksendal. *Stochastic Differential Equations: An Introduction with Applications (Universitext)*. Springer, 6th edition, January 2014.
- [6] Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations, 2021.
- [7] Simo Särkkä and Arno Solin. *Applied Stochastic Differential Equations*. Institute of Mathematical Statistics Textbooks. Cambridge University Press, 2019.

A Conservation Laws

A.1 Flow Map

Theorem 29 *Two important theorems in calculus:*

1. Divergence Theorem:

$$\int_{\Omega} \nabla \cdot \mathbf{F} dx = \int_{\partial\Omega} \mathbf{F} \cdot \mathbf{n} dS \quad (127)$$

2. Reynolds Transport Theorem:

$$\frac{d}{dt} \int_{\Omega(t)} f(t, x) dx = \int_{\Omega(t)} \frac{\partial f}{\partial t} dx + \int_{\partial\Omega(t)} f(t, x) \mathbf{v} \cdot \mathbf{n} dS \quad (128)$$

where \mathbf{v} is the velocity at $\partial\Omega(t)$.

Here the $\Omega(t)$ is the domain of the flow, and the $\partial\Omega(t)$ is the boundary of the flow, which is described by the flow map ϕ_s^t . Here is the definition.

Definition 16 (Flow Map) *Assume a description of some characteristic of particle \mathbf{P} , like the position or the boundary, as $\mathbf{x} \in \mathcal{R}^m$, then we have a flow map $\phi_s^t(\mathbf{x}) \in \mathcal{R}^m$, which means that the flow transmits the characteristic(position) \mathbf{x} from \mathbf{x} at s to $\phi_s^t(\mathbf{x})$ at t , controlled by the vector field(velocity field) $\mathbf{F} : \mathcal{R}^m \times \mathcal{R} \rightarrow \mathcal{R}^m$:*

$$\begin{cases} \frac{d\phi_s^t(\mathbf{x})}{dt} = \mathbf{F}(\phi_s^t(\mathbf{x}), t) \\ \phi_s^s(\mathbf{x}) = \mathbf{x} \end{cases} \quad (129)$$

A.1.1 Conservation Laws

If we assume $\Omega(t)$ is composed of particles, i.e. $\Omega(t) = \phi_{t_0}^t(\Omega)$ (when $t = t_0$, $\Omega(t_0) = \Omega$), then we by **conservation of mass**, we have the following theorem:

Theorem 30 (Continuity Equation) *By conservation of mass, i.e. $\int_{\Omega(t)} \rho(t, \mathbf{x}) d\mathbf{x} = C$, we have:*

$$\begin{aligned} \frac{d}{dt} \int_{\Omega(t)} \rho(t, \mathbf{x}) d\mathbf{x} &= \int_{\Omega(t)} \frac{\partial \rho}{\partial t} d\mathbf{x} + \int_{\partial\Omega(t)} \rho(t, \mathbf{x}) \mathbf{u} \cdot \mathbf{n} dS \\ &= \int_{\Omega(t)} \left(\frac{\partial \rho}{\partial t} + \nabla \cdot (\rho \mathbf{u}) \right) d\mathbf{x} = 0 \end{aligned} \quad (130)$$

Therefore:

$$\frac{\partial \rho}{\partial t} + \nabla \cdot (\rho \mathbf{u}) = 0 \quad (131)$$

which is also called **continuity equation**.

Theorem 31 (Conservation of Momentum) *By conservation of momentum, i.e.*

$$\frac{d}{dt} \int_{\Omega(t)} \rho(t, \mathbf{x}) \mathbf{v}(t, \mathbf{x}) d\mathbf{x} = - \int_{\partial\Omega(t)} p \cdot \mathbf{n} dS \quad (132)$$

we have:

$$\frac{\partial(\rho \mathbf{v})}{\partial t} + \nabla \cdot (\rho \mathbf{v} \otimes \mathbf{v} + p) = 0 \quad (133)$$

where p is the pressure.

Theorem 32 (Conservation of Energy)

$$\frac{\partial E}{\partial t} + \nabla \cdot (\mathbf{v}(E + p)) = 0 \quad (134)$$

Then we have can get the Euler's equation:

Theorem 33 (Euler's Equation) *The Euler's equation is given by:*

$$\frac{\partial}{\partial t} \begin{bmatrix} \rho \\ \rho \mathbf{v} \\ E \end{bmatrix} + \nabla \cdot \begin{bmatrix} \rho \mathbf{v} \\ \rho \mathbf{v} \otimes \mathbf{v} + p \\ \mathbf{v}(E + p) \end{bmatrix} = 0 \quad (135)$$

So the general form of conservation laws is given by: suppose $U \in \mathcal{R}^d$ is the conserved quantity, F is $\mathcal{R}^d \rightarrow \mathcal{R}^d$ is the flux, then we have:

$$\frac{\partial U}{\partial t} + \nabla \cdot (F(U)) = 0 \quad (136)$$

B Priori

B.1 Hilbert space-valued random variable

Definition 17 ($L^p(\Omega, H)$ space) Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space and H is a Hilbert space with norm $\|\cdot\|$. Then $\mathcal{L}^p(\Omega, H)$ with $1 \leq p < \infty$ is the space of H -valued \mathcal{F} -measurable random variables $X : \Omega \rightarrow H$ with $\mathbf{E}[\|X\|^p] < \infty$ and a Banach space with norm:

$$\|X\|_{\mathcal{L}^p(\Omega, H)} := \left(\int_{\Omega} \|X(\omega)\|^p dP(\omega) \right)^{\frac{1}{p}} = \mathbf{E}[\|X\|^p]^{\frac{1}{p}} \quad (137)$$

Then we can define the inner product:

$$\langle X, Y \rangle_{\mathcal{L}^2(\Omega, H)} := \int_{\Omega} \langle X(\omega), Y(\omega) \rangle dP(\omega) \quad (138)$$

Definition 18 (uncorrelated, covariance operator) Let H be a Hilbert space. A linear operator $\mathcal{C} : H \rightarrow H$ is the covariance of H -valued random variables X and Y if

$$\langle \mathcal{C}\phi, \psi \rangle = \text{Cov}(\langle X, \phi \rangle, \langle Y, \psi \rangle), \forall \phi, \psi \in H \quad (139)$$

specially, we show that in finite dimensional case, the covariance matrix coincides with the covariance operator. when $H = \mathbb{R}^d$,

$$\begin{aligned} \text{Cov}(\langle X, \phi \rangle, \langle Y, \psi \rangle) &= \text{Cov}(\phi^T X, \psi^T Y) \\ &= \mathbf{E}[\phi^T (X - \mu_X)(Y - \mu)^T \psi] = \phi^T \mathbf{E}[(X - \mu_X)(Y - \mu)^T] \psi \\ &= \langle \mathcal{C}\phi, \psi \rangle \end{aligned} \quad (140)$$

Definition 19 (H -valued Gaussian random variable) Let H be a Hilbert space. An H -valued random variable X is Gaussian if $\langle X, \phi \rangle$ is a real-valued Gaussian random variable for all $\phi \in H$.

B.2 Hilbert-Schmidt operator

Definition 20 (Hilbert-Schmidt operator) Let U, H be two separable Hilbert spaces with norms $\|\cdot\|, \|\cdot\|_U$ respectively. For an orthonormal basis $\{\phi_j\}$ of U , define the Hilbert-Schmidt norm:

$$\|L\|_{\text{HS}(U, H)} := \left(\sum_{j=1}^{\infty} \|L\phi_j\|_H^2 \right)^{\frac{1}{2}} \quad (141)$$

where $\text{HS}(U, H) := \{L \in \mathcal{L}(U, H) : \|L\|_{\text{HS}(U, H)} < \infty\}$ is a Banach space with Hilbert-Schmidt norm. And $L \in \text{HS}(U, H)$ is called Hilbert-Schmidt operator.

Definition 21 (Integral operator with kernel G) For a domain D and a kernel $G \in L^2(D \times D)$, define the integral operator L by

$$(Lu)(x) = \int_D G(x, y)u(y)dy, x \in D, u \in L^2(D) \quad (142)$$

Furthermore, L is a Hilbert-Schmidt operator.

B.3 Operator theory

Theorem 34 (Sobolev embedding theorem) 1. Let $W^{r,p}(\mathbf{R}^n)$. Here k is a non-negative integer and $1 \leq p < \infty$. If $k > \ell, p < n$ and $1 \leq p < q < \infty$ are two real numbers such that $\frac{1}{p} - \frac{r}{n} = \frac{1}{q} - \frac{\ell}{n}$, then

$$W^{r,p}(\mathbf{R}^n) \subseteq W^{\ell,q}(\mathbf{R}^n) \quad (143)$$

Specially, if $\ell = 0$, then $\frac{1}{p} - \frac{r}{n} = \frac{1}{q}$, then $W^{r,p}(\mathbf{R}^n) \subseteq L^q(\mathbf{R}^n)$.

2. If $n < pr$ and $\frac{1}{p} - \frac{r}{n} = -\frac{s+\alpha}{n}$, then $W^{r,p}(\mathbf{R}^n) \subseteq C^{s,\alpha}(\mathbf{R}^n)$.

Definition 22 (domain of operator) For a linear operator $A : \mathcal{D}(A) \subset H \rightarrow H$, the domain of A is defined as $\mathcal{D}(A)$

Theorem 35 (Dirichlet Boundary Condition) Consider the Dirichlet problem for Poisson equation: for $f \in L^2(0, 1)$, find $u \in H^2(0, 1)$ s.t.

$$\begin{aligned} u_{xx} &= f, \quad x \in (0, 1) \\ u(0) &= u(1) = 0 \end{aligned} \quad (144)$$

We also assume $u \in H_0^1(0, 1)$. By Sobolev embedding theorem, $u \in H_0^1(0, 1) \subset C([0, 1])$. Then, Laplacian with Dirichlet conditions can be defined as:

$$Au := -u_{xx}, u \in \mathcal{D}(A) = H^2(0, 1) \cap H_0^1(0, 1) \quad (145)$$

Definition 23 (Periodic Boundary Condition) ...

Definition 24 If A is a linear operator from $\mathcal{D}(A) \subset H$ to Hilbert space H , with an orthonormal basis of eigenfunctions $\{\phi_j\}$ and corresponding increasing eigenvalues $\{\lambda_j\}$, then A^α is defined as:

$$A^\alpha u = \sum_{j=1}^{\infty} \lambda_j^\alpha \langle u, \phi_j \rangle \phi_j \quad (146)$$

and the domain $\mathcal{D}(A^\alpha)$ is the set of all $u \in H$ such that $A^\alpha u \in H$.